



# Learning from Imbalanced Data An Application to Bank Fraud Detection

Guillaume Metzler

Univ. Lyon, UJM-Saint-Etienne, CNRS, Institut d'Optique Graduate School, Laboratoire  
Hubert Curien UMR 5516, F-42023, SAINT-ETIENNE, France  
Blitz Business Services, VILLEFONTAINE, France

25 September 2019

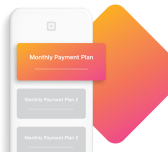
Marianne CLAUSEL	Professeure, Université de Lorraine	Rapporteuse
Marc TOMMASI	Professeur, Université de Lille	Rapporteur
Yves GRANDVALET	Professeur, Université de Technologie de Compiègne	Examineur
Élisa FROMONT	Professeure, Université de Rennes I	Co-encadrante
Amaury HABRARD	Professeur, Université de Saint-Étienne	Co-encadrant
Marc SEBBAN	Professeur, Université de Saint-Étienne	Directeur
Xavier BADICHE	Président Directeur Général de Blitz Business Services	Blitz BS
Brahim BELKASMI	Ingénieur R&D, Blitz Business Services	Blitz BS

# Context of Thesis

Blitz company

## Blitz activities

Buy now.  
Pay in monthly  
installments.



Payment facilities



Smooth checkout flow



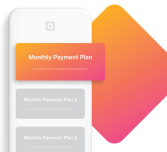
Securing cheque  
transactions

# Context of Thesis

Blitz company

## Blitz activities

Buy now.  
Pay in monthly  
installments.



Payment facilities



Smooth checkout flow



Securing cheque  
transactions

## Other activities:

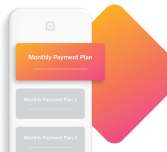
- Assistance with PV management
- Assistance in staff management

# Context of Thesis

Blitz company

## Blitz activities

Buy now.  
Pay in monthly  
installments.



Payment facilities



Smooth checkout flow



Securing cheque  
transactions

## Other activities:

- Assistance with PV management
- Assistance in staff management

→ This work Focus on the topic of securing cheque transactions ...

# Context of Thesis

## Check fraud detection

### What is cheque fraud ?



- Unpaid cheque (no money on bank account)
- False cheque
  - Not the real identity
  - Incorrect number series in the CMC7

# Context of Thesis

## Check fraud detection

### What is cheque fraud ?



- Unpaid cheque (no money on bank account)
- False cheque
  - Not the real identity
  - Incorrect number series in the CMC7

### Some statistics :

10 months of transactions  
(03/20/2016 to 10/21/2016)

- around 3.2 millions of transactions
- for 195 millions of euros
- 20 000 are frauds or unpaid (0.6%)
- represent 2 millions of euros (1.1%)

# Context of Thesis

## Check fraud detection

### What is cheque fraud ?



- Unpaid cheque (no money on bank account)

- False cheque

Not the real identity

Incorrect number series in the CMC7

### Some statistics :

10 months of transactions  
(03/20/2016 to 10/21/2016)

- around 3.2 millions of transactions
- for 195 millions of euros
- 20 000 are frauds or unpaid (0.6%)
- represent 2 millions of euros (1.1%)

... more precisely on the topic of learning from imbalanced data

1. Introduction on Learning From Imbalanced Data
2. A Geometrical Approach based on the Distance to Positives
  - 2.1 Building Risky Areas

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*
  - 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*
3. An Approach based on Cost-Sensitive Learning
  - 3.1 Optimizing F-measure by Cost-Sensitive Classification

*CONE: "From Cost-Sensitive Classification to Tight F-Measure Bounds"*
  - 3.2 Improving the Benefits of Mass Distribution

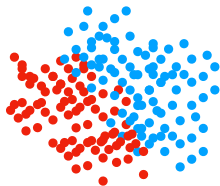
*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*
4. Conclusion and Perspectives



# Learning from Imbalanced Data

## Balanced vs. Imbalanced

**Balanced dataset**

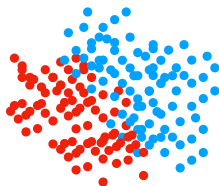


Positives  $\simeq$  Negatives

# Learning from Imbalanced Data

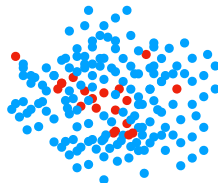
## Balanced vs. Imbalanced

Balanced dataset



Positives  $\simeq$  Negatives

Imbalanced dataset



Positives  $\ll$  Negatives

Minimizing a surrogate of  $\frac{1}{m} \sum_{i=1}^m 1_{\{\hat{y}_i \neq y_i\}}$  leads to:

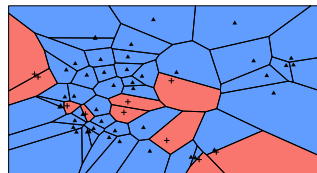
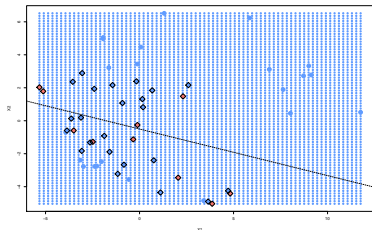
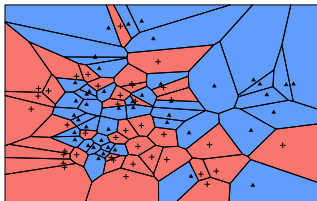
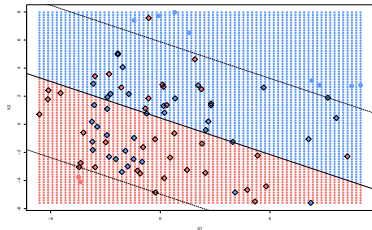
focus on both classes

focus on majority examples

# Learning from Imbalanced Data

## Impact of Imbalance

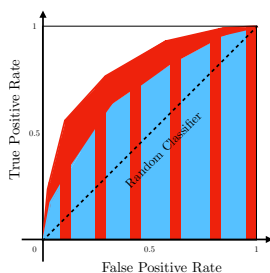
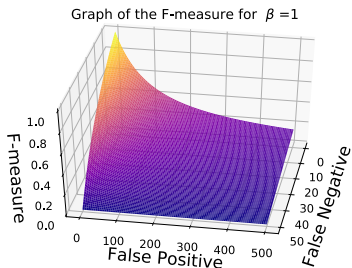
Example of linear SVM and  $k$ -NN with 50% and 20% of positives.



# Learning from Imbalanced Data

## Performance Measures

### Use appropriate measures



$$F_{\beta} = \frac{(1 + \beta^2)(P - FN)}{(1 + \beta^2)P - FN + FP}$$

$$\mathbb{P}[f(x_+) > f(x_-)]$$

**G-mean**

**Precision**

**MCC**

**Mean Average Precision**

**False Positive Rate**

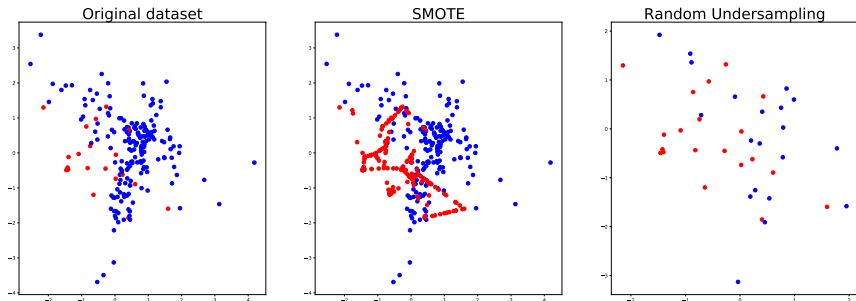
**Recall**

**Average Precision**

# Learning from Imbalanced Data

Balance the two classes

## Use sampling strategies



- Oversampling: Random - SMOTE - BorderSMOTE, ...
- Undersampling: Random - Tomek Link - ENN, ...

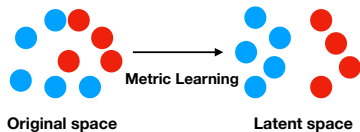
# Learning from Imbalanced Data

## Representation and Cost-Sensitive Learning

### Distance and representation

$$d_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')},$$

where  $\mathbf{M}$  is **PSD**.



### Cost-sensitive learning

$$C_{TP} = 0 \quad C_{FN} = c$$

$$C_{FP} = 1 - c \quad C_{TN} = 0$$

$\rightarrow c \simeq 1$  to encourage low miss-classification on positives.

$$\ell(y, h(\mathbf{x})) = c \cdot y \cdot (1 - h(\mathbf{x})) + (1 - c) \cdot (1 - y) \cdot h(\mathbf{x})$$

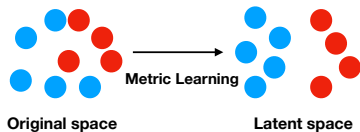
# Learning from Imbalanced Data

## Representation and Cost-Sensitive Learning

### Distance and representation

$$d_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')},$$

where  $\mathbf{M}$  is **PSD**.



### Cost-sensitive learning

$$C_{TP} = 0 \quad C_{FN} = c$$

$$C_{FP} = 1 - c \quad C_{TN} = 0$$

→  $c \simeq 1$  to encourage low miss-classification on positives.

$$\ell(y, h(\mathbf{x})) = c \cdot y \cdot (1 - h(\mathbf{x})) + (1 - c) \cdot (1 - y) \cdot h(\mathbf{x})$$

→ *Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guarantees*, PRL, 2018.

→ *An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data*, IC-TAI, 2019.

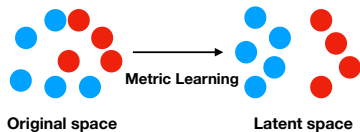
# Learning from Imbalanced Data

## Representation and Cost-Sensitive Learning

### Distance and representation

$$d_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')},$$

where  $\mathbf{M}$  is **PSD**.



### Cost-sensitive learning

$$C_{TP} = 0 \quad C_{FN} = c$$

$$C_{FP} = 1 - c \quad C_{TN} = 0$$

→  $c \simeq 1$  to encourage low misclassification on positives.

$$\ell(y, h(\mathbf{x})) = c \cdot y \cdot (1 - h(\mathbf{x})) + (1 - c) \cdot (1 - y) \cdot h(\mathbf{x})$$

→ *Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guarantees*, PRL, 2018.

→ *An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data*, IC-TAI, 2019.

→ *From Cost-Sensitive Classification to Tight F-Measure Bounds*, AISTATS, 2019.

→ *Tree-based Cost Sensitive Methods for Fraud Detection in Imbalanced Data*, IDA, 2018.



1. Introduction on Learning From Imbalanced Data
- 2. A Geometrical Approach based on the Distance to Positives**
  - 2.1 Building Risky Areas**

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*
  - 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*
3. An Approach based on Cost-Sensitive Learning
  - 3.1 Optimizing F-measure by Cost-Sensitive Classification

*CONE: "From Cost-Sensitive Classification to Tight F-Measure Bounds"*
  - 3.2 Improving the Benefits of Mass Distribution

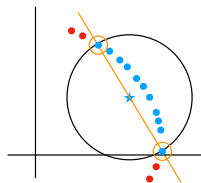
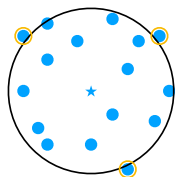
*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*
4. Conclusion and Perspectives

# $ME^2$ : Learning Risky Areas

## Hypothesis

Frauds are close to each other, they form small groups in the feature space

Given a set of  $m$  unlabelled points, find the center  $\mathbf{c}$  and the **smallest** radius  $R$  of the ball that includes the data (Tax and Duin, 2004).



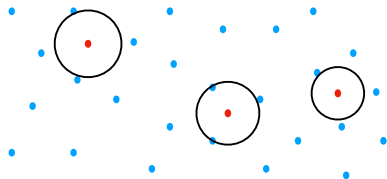
$$\begin{aligned} \min_{\mathbf{c}, R, \xi} \quad & R^2 + \frac{\mu}{m} \sum_{i=1}^m \xi_i, \\ \text{s.t.} \quad & \|\mathbf{x}_i - \mathbf{c}\|_2^2 \leq R^2 + \xi_i, \quad \forall i, \\ & \xi_i \geq 0 \quad \forall i. \end{aligned}$$

$$\begin{aligned} \min_{\mathbf{c}, \rho, \xi} \quad & \frac{1}{2} \|\mathbf{c}\|_2^2 + \frac{1}{\nu m} \sum_{i=1}^m \xi_i - \rho - \frac{1}{2} \|\mathbf{x}_i\|_2^2, \\ \text{s.t.} \quad & \mathbf{c}^T \mathbf{x}_i \geq \rho + \frac{1}{2} \|\mathbf{x}_i\|_2^2, \\ & \xi_i \geq 0 \quad \forall i. \end{aligned}$$

Being in the ball  $\iff$  being above the hyperplane

# $ME^2$ : Learning Risky Areas

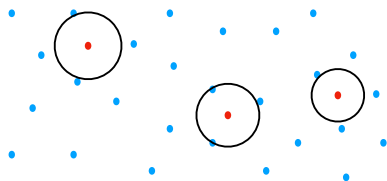
From  $MIB$  to  $ME^2$



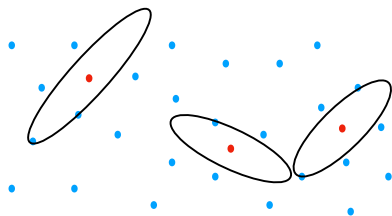
- Use the idea of MIB to create MEB
- One model per positive instance
- Require few positive neighbors

# $ME^2$ : Learning Risky Areas

From  $MIB$  to  $ME^2$



↓ Learning a Metric



- Use the idea of MIB to create MEB
- One model per positive instance
- Require few positive neighbors

- From balls to ellipsoids
  - Increase decision boundary
- Maximum Excluding Ellipsoids

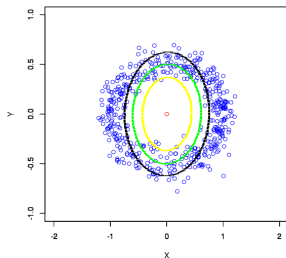
# $ME^2$ : Learning Risky Areas

Optimization problem

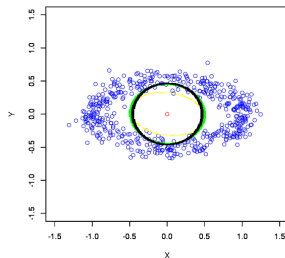
$$\begin{aligned} \min_{R, \mathbf{M}, \xi} \quad & \frac{1}{m} \sum_{i=1}^m \xi_i + \mu (B - R)^2 + \lambda \|\mathbf{M} - \mathbf{I}\|_{\mathcal{F}^2}, \\ \text{s.t.} \quad & \|\mathbf{x}_i - \mathbf{c}\|_{\mathbf{M}}^2 \geq R - \xi_i, \quad \forall i = 1, \dots, m, \\ & \xi_i \geq 0, \quad \forall i = 1, \dots, m \\ & 0 \leq R \leq B, \end{aligned}$$

error terms (in terms of distances)

regularization term



Influence of  $\mu$



Influence of  $\lambda$

# $ME^2$ : Learning Risky Areas

## Dual formulation

- express the Lagrangian  $\mathcal{L}$  including the constraints
- expression of primal variables w.r.t. dual ones:
  1. derivative of  $\mathcal{L}$  w.r.t. primal variables
  2. set derivatives to 0

# $ME^2$ : Learning Risky Areas

## Dual formulation

- express the Lagrangian  $\mathcal{L}$  including the constraints
- expression of primal variables w.r.t. dual ones:

1. derivative of  $\mathcal{L}$  w.r.t. primal variables
2. set derivatives to 0

One of these derivatives gives:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{M}} = 0 \implies \mathbf{M} = \mathbf{I} + \frac{1}{2\lambda} \sum_{i=1}^m \alpha_k (\mathbf{x}_k - \mathbf{c})(\mathbf{x}_k - \mathbf{c})^T.$$

→ **M is Positive Semi Definite for free**

# $ME^2$ : Learning Risky Areas

## Theoretical Guarantees

Using stability framework (Bousquet and Elisseeff, 2002)

$$\mathcal{R}(\mathbf{M}, R) \leq \mathcal{R}_S(\mathbf{M}, R) + \mathcal{O} \left( \frac{1}{\min(\mu, \lambda)} \sqrt{\frac{\ln(1/\delta)}{2m}} \right),$$

where  $\mathcal{R}_S(\mathbf{M}, R) = \frac{1}{m} \sum_{i=1}^m [R - \|\xi - \mathbf{c}\|_{\mathbf{M}}^2]_+$ .



# $ME^2$ : Learning Risky Areas

## Theoretical Guarantees

Using stability framework (Bousquet and Elisseeff, 2002)

$$\mathcal{R}(\mathbf{M}, R) \leq \mathcal{R}_S(\mathbf{M}, R) + \mathcal{O}\left(\frac{1}{\min(\mu, \lambda)} \sqrt{\frac{\ln(1/\delta)}{2m}}\right),$$

where  $\mathcal{R}_S(\mathbf{M}, R) = \frac{1}{m} \sum_{i=1}^m [R - \|\xi - \mathbf{c}\|_{\mathbf{M}}^2]_+$ .

the true risk on the underlying and unknown distribution

# $ME^2$ : Learning Risky Areas

## Theoretical Guarantees

Using stability framework (Bousquet and Elisseeff, 2002)

$$\mathcal{R}(\mathbf{M}, R) \leq \mathcal{R}_S(\mathbf{M}, R) + \mathcal{O}\left(\frac{1}{\min(\mu, \lambda)} \sqrt{\frac{\ln(1/\delta)}{2m}}\right),$$

where  $\mathcal{R}_S(\mathbf{M}, R) = \frac{1}{m} \sum_{i=1}^m [R - \|\xi - \mathbf{c}\|_{\mathbf{M}}^2]_+$ .

the true risk on the underlying and unknown distribution

the empirical risk over the sample  $S$

# $ME^2$ : Learning Risky Areas

## Theoretical Guarantees

Using stability framework (Bousquet and Elisseeff, 2002)

$$\mathcal{R}(\mathbf{M}, R) \leq \mathcal{R}_S(\mathbf{M}, R) + \mathcal{O} \left( \frac{1}{\min(\mu, \lambda)} \sqrt{\frac{\ln(1/\delta)}{2m}} \right),$$

where  $\mathcal{R}_S(\mathbf{M}, R) = \frac{1}{m} \sum_{i=1}^m [R - \|\xi - \mathbf{c}\|_{\mathbf{M}}^2]_+$ .

the true risk on the underlying and unknown distribution

the empirical risk over the sample  $S$

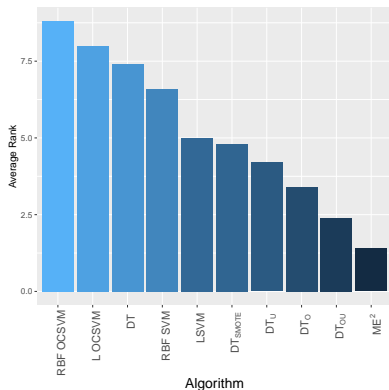
generalization gap of the learned model: depends on the complexity of the model

# ME<sup>2</sup>: Learning Risky Areas

## Experimental Results

Comparison with standards algorithms on imbalanced datasets

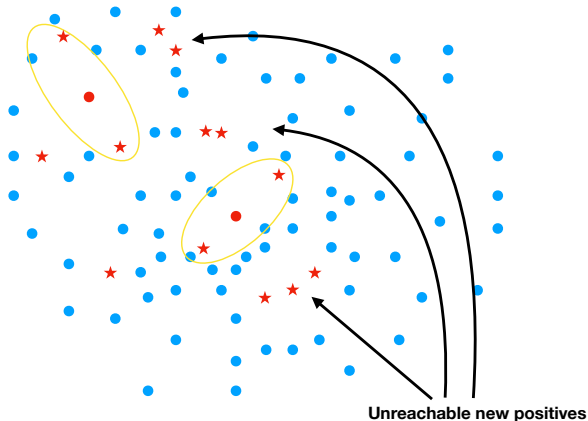
Dataset	Nb. of ex.	% Pos.
Wine	1 599	3.3
Abalone17	2 338	2.5
Yeast6	1 484	2.4
Abalone20	1 916	1.4
Blitz	15 000	1.0



**Lower Rank:** able to reach better performance

# $ME^2$ : Learning Risky Areas

Limitation of  $ME^2$



Find a way to increase the influence zone of positives

## 1. Introduction on Learning From Imbalanced Data

## 2. A Geometrical Approach based on the Distance to Positives

### 2.1 Building Risky Areas

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*

### 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*

## 3. An Approach based on Cost-Sensitive Learning

### 3.1 Optimizing F-measure by Cost-Sensitive Classification

*CONE: "From Cost-Sensitive Classification to Tight F-Measure Bounds"*

### 3.2 Improving the Benefits of Mass Distribution

*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*

## 4. Conclusion and Perspectives

# $\gamma$ - $k$ -NN : a revisit of the $k$ -NN

## Presentation of $\gamma$ - $k$ -NN

### Observations

Imbalanced setting  $\rightarrow$  low density of positives

low density of positives  $\rightarrow$  small influence area

# $\gamma$ - $k$ -NN : a revisit of the $k$ -NN

## Presentation of $\gamma$ - $k$ -NN

### Observations

Imbalanced setting  $\rightarrow$  low density of positives

low density of positives  $\rightarrow$  small influence area

### Idea

Bring points closer to positives by modifying their distances



# $\gamma$ -k-NN : a revisit of the $k$ -NN

## Presentation of $\gamma$ -k-NN

### Observations

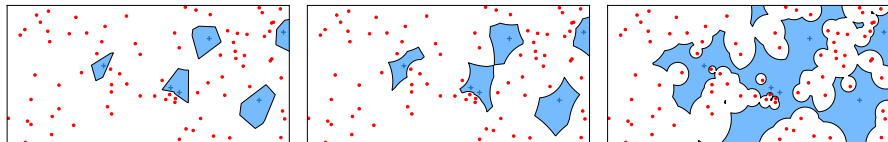
Imbalanced setting  $\rightarrow$  low density of positives

low density of positives  $\rightarrow$  small influence area

### Idea

Bring points closer to positives by modifying their distances

$$d_{\gamma}(\mathbf{x}, \mathbf{x}_i) = \begin{cases} d(\mathbf{x}, \mathbf{x}_i) & \text{if } y_i = -1, \\ \gamma \cdot d(\mathbf{x}, \mathbf{x}_i) & \text{if } y_i = +1. \end{cases}$$



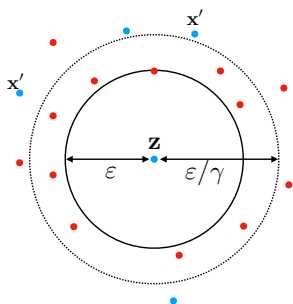
# $\gamma$ -k-NN: a revisit of the $k$ -NN

## Study of $\gamma$ parameter

### Importance of the parameter $\gamma$

Probability of False Negative

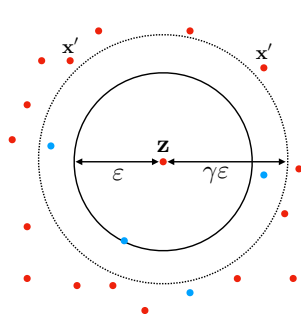
$$\gamma \leq 1$$



$$FN_{\gamma}(\mathbf{z}) = \left(1 - \mathbb{P}(\mathbf{x}' \in \mathcal{S}_{\frac{\epsilon}{\gamma}})\right)^{m_+}$$

Probability of False Positive

$$\gamma \geq 1$$



$$FP_{\gamma}(\mathbf{z}) = \left(1 - \mathbb{P}(\mathbf{x}' \in \mathcal{S}_{\epsilon\gamma})\right)^{m_-}$$

**Choose  $\gamma < 1$  in Imbalanced settings**

## $\gamma$ -k-NN: a revisit of the $k$ -NN

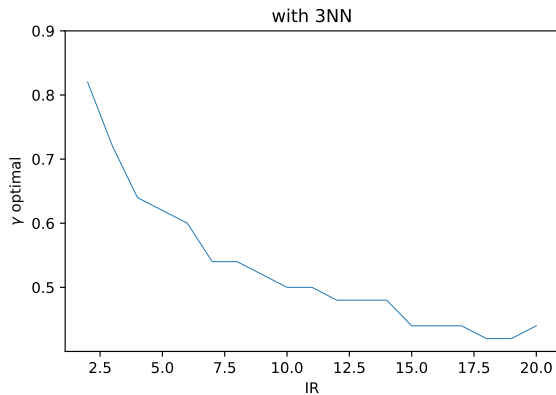
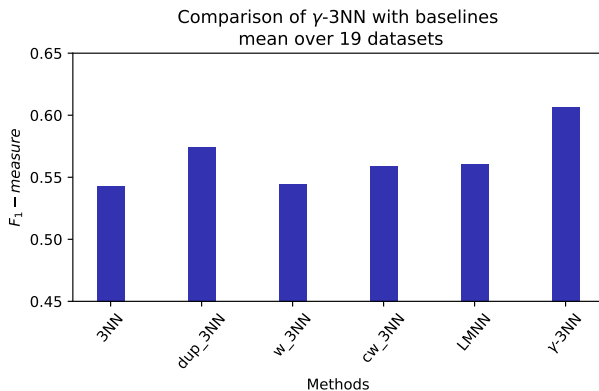


Illustration of the optimal  $\gamma$  with respect to the IR on Balance dataset

# $\gamma$ -k-NN: a revisit of the $k$ -NN

## Experimental results 1/2

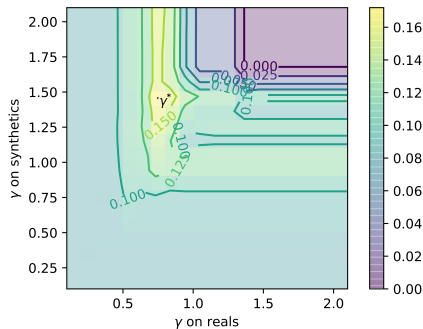


→ perform better and even better than a Metric Learning approach.

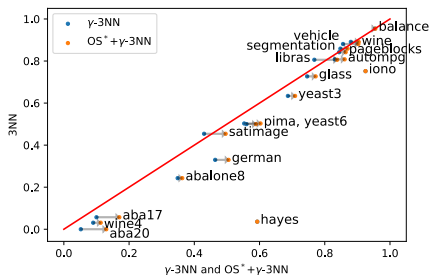
# $\gamma$ -k-NN: a revisit of the $k$ -NN

Experimental results 2/2

## Behaviour of $\gamma$ -k-NN combined with an over-sampler.



reals:  $\gamma < 1$   
synthetics:  $\gamma > 1$



Coupling with sampling strategies  
improves the algorithm

1. Introduction on Learning From Imbalanced Data
2. A Geometrical Approach based on the Distance to Positives
  - 2.1 Building Risky Areas

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*
  - 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*
3. An Approach based on Cost-Sensitive Learning
  - 3.1 Optimizing F-measure by Cost-Sensitive Classification

**CONE**: *"From Cost-Sensitive Classification to Tight F-Measure Bounds"*
  - 3.2 Improving the Benefits of Mass Distribution

*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*
4. Conclusion and Perspectives

# CONE: an Algorithm for F-measure Optimization

## F-measure

**Objective:** find a way to optimize the F-measure  $F_\beta$

$$F_\beta = \frac{(1 + \beta^2)(P - FN)}{(1 + \beta^2)P - FN + FP} = \frac{(1 + \beta^2)(P - e_1)}{(1 + \beta^2)P - e_1 + e_2}.$$

Two important quantities:  $e_1 = FN$  et  $e_2 = FP$  linked to the empirical risk  $\mathcal{R}$ .

# CONE: an Algorithm for F-measure Optimization

## F-measure

**Objective:** find a way to optimize the F-measure  $F_\beta$

$$F_\beta = \frac{(1 + \beta^2)(P - FN)}{(1 + \beta^2)P - FN + FP} = \frac{(1 + \beta^2)(P - e_1)}{(1 + \beta^2)P - e_1 + e_2}.$$

Two important quantities:  $e_1 = FN$  et  $e_2 = FP$  linked to the empirical risk  $\mathcal{R}$ .

How to make the link between  $F_\beta$  and  $\mathcal{R}$ ?



# CONE: an Algorithm for F-measure Optimization

## F-measure

**Objective:** find a way to optimize the F-measure  $F_\beta$

$$F_\beta = \frac{(1 + \beta^2)(P - FN)}{(1 + \beta^2)P - FN + FP} = \frac{(1 + \beta^2)(P - e_1)}{(1 + \beta^2)P - e_1 + e_2}.$$

Two important quantities:  $e_1 = FN$  et  $e_2 = FP$  linked to the empirical risk  $\mathcal{R}$ .

How to make the link between  $F_\beta$  and  $\mathcal{R}$ ?

→ **Pseudo linearity of the F-measure !**

# CONE: an Algorithm for F-measure Optimization

## Related work

- Based on previous work published in 2014 at NIPS (Parambath et al., 2014)
- Use the pseudo-linearity of the F-measure
- Derive bounds on optimality of  $F_\beta$
- Algorithmic : grid approach

# CONE: an Algorithm for F-measure Optimization

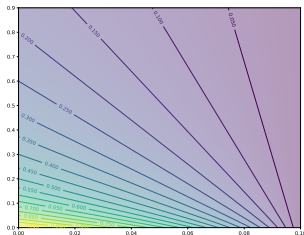
## Related work

- Based on previous work published in 2014 at NIPS (Parambath et al., 2014)
- Use the pseudo-linearity of the F-measure
- Derive bounds on optimality of  $F_\beta$
- Algorithmic : grid approach

→ **Extend the existing work from both theoretical and practical aspect**

# CONE: an Algorithm for F-measure Optimization

A pseudo linear function



- $F_\beta$  level sets are hyperplanes in the  $(e_1, e_2)$ -space:

$$\forall t \in [0, 1], F_\beta(\mathbf{e}) = t \iff \exists \mathbf{a}, b \text{ t.q. } \langle \mathbf{a}(t), \mathbf{e} \rangle + b(t) = 0.$$

- $\mathbf{a}$  : weights assigned to the errors
- $\langle \mathbf{a}(t), \mathbf{e} \rangle$  : weighted version of  $\mathcal{R}$ .

→ Good choice of  $t \iff$  Optimizing  $F_\beta$ .

# CONE: an Algorithm for F-measure Optimization

Deriving a bound 1/2

- Write the difference of F-measures between  $\mathbf{e}$  and  $\mathbf{e}'$

$$F(\mathbf{e}') - F(\mathbf{e}) = \Phi_{\mathbf{e}} \cdot \langle \mathbf{a}(F(\mathbf{e}')), \mathbf{e} - \mathbf{e}' \rangle, \quad \Phi_{\mathbf{e}} = \frac{1}{(1 + \beta^2)P - e_1 + e_2}.$$

- Bound this difference using:
  1. linearity of the inner product
  2. sub-optimality  $\varepsilon_1$  of the learned hypothesis

$$F(\mathbf{e}') - F(\mathbf{e}) \leq \Phi_{\mathbf{e}} \varepsilon_1 + \Phi_{\mathbf{e}} \cdot (e_2 - e_1 - (e'_2 - e'_1)) (t' - t).$$

**Problem:**  $\mathbf{e}(t') = \mathbf{e}' = (e'_1, e'_2)$  is unknown

# CONE: an Algorithm for F-measure Optimization

## Deriving a bound 2/2

→ Bound the difference  $e'_2 - e'_1$

- When  $t' < t$ :

$$M_{max} = \max_{\mathbf{e}'' \in \mathcal{E}(\mathcal{H})} (e''_2 - e''_1) \\ \text{s.t. } F(\mathbf{e}'') > F(\mathbf{e})$$

$$F(\mathbf{e}') \leq F(\mathbf{e}) + \Phi_{\mathbf{e}} \varepsilon_1 + \Phi_{\mathbf{e}} \cdot (e_2 - e_1 - M_{max})(t' - t),$$

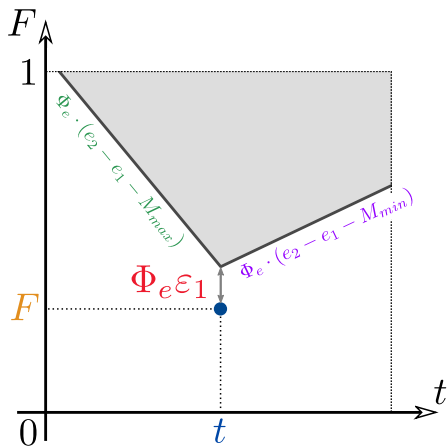
- When  $t' > t$ :

$$M_{min} = \min_{\mathbf{e}'' \in \mathcal{E}(\mathcal{H})} (e''_2 - e''_1) \\ \text{s.t. } F(\mathbf{e}'') > F(\mathbf{e})$$

$$F(\mathbf{e}') \leq F(\mathbf{e}) + \Phi_{\mathbf{e}} \varepsilon_1 + \Phi_{\mathbf{e}} \cdot (e_2 - e_1 - M_{min})(t' - t),$$

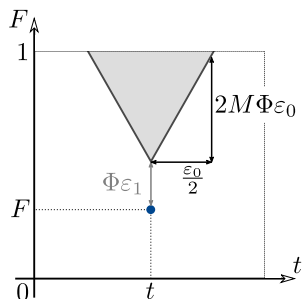
# CONE: an Algorithm for F-measure Optimization

An asymmetric cone



# CONE : an Algorithm for F-measure Optimization

Existing results



Interpretation existing bound  
of Parambath et al. (2014)

$$F(\mathbf{e}') \leq F(\mathbf{e}) + \Phi \cdot (2\varepsilon_0 M + \varepsilon_1)$$

$$F(\mathbf{e}') \leq F(\mathbf{e}) + \Phi\varepsilon_1 + 4M\Phi|t' - t|.$$

où

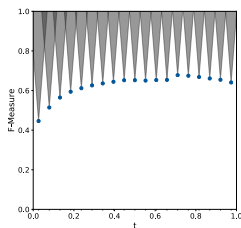
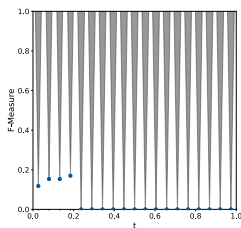
- $\varepsilon_0$ : distance to optimal weights
- $M = \max_{\mathbf{e}''} \|\mathbf{e}''\|_2$
- $\Phi = (\beta^2 P)^{-1}$  independent from  $\mathbf{e}$



# CONE: an Algorithm for F-measure Optimization

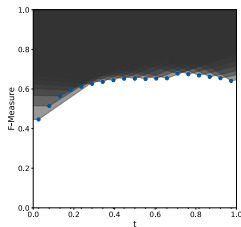
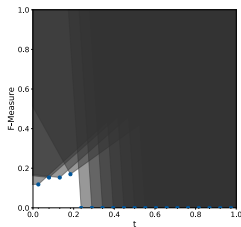
Bounds comparison: Parambath et al. (2014) vs Our

Similar bounds ...



Abalone 12

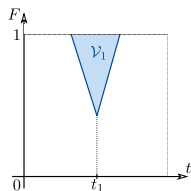
Adult



... with highly different slopes.

# CONE: an Algorithm for F-measure Optimization

An iterative algorithm

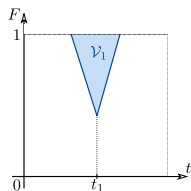


Step 1

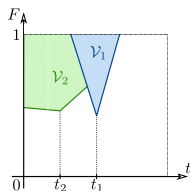
**Step 1:** Take the middle of the  $t$ -space of reasearch:  $t_1 = 0.5$   
→ Highest values of  $F$  in  $[0, t_1]$

# CONE: an Algorithm for F-measure Optimization

An iterative algorithm



Step 1



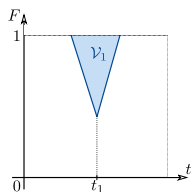
Step 2

**Step 1:** Take the middle of the  $t$ -space of reasearch:  $t_1 = 0.5$   
→ Highest values of  $F$  in  $[0, t_1]$

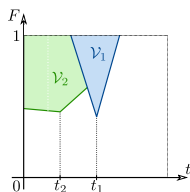
**Step 2:** Choose  $t_2$  in the middle of  $[0, t_1]$   
→ Highest values of  $F$  in  $[t_1, 1]$

# CONE: an Algorithm for F-measure Optimization

An iterative algorithm



Step 1

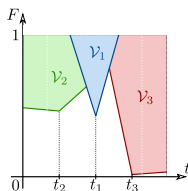


Step 2

**Step 1:** Take the middle of the  $t$ -space of reasearch:  $t_1 = 0.5$   
→ Highest values of  $F$  in  $[0, t_1]$

**Step 2:** Choose  $t_2$  in the middle of  $[0, t_1]$   
→ Highest values of  $F$  in  $[t_1, 1]$

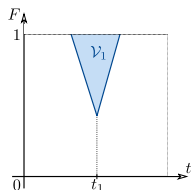
**Step 3:** Choose  $t_3$  in the middle of  $[t_1, 1]$   
→ Highest values of  $F$  in  $[t_1, t_3]$



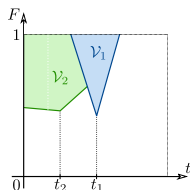
Step 3

# CONE: an Algorithm for F-measure Optimization

An iterative algorithm



Step 1



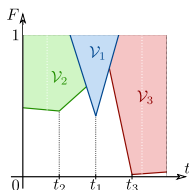
Step 2

**Step 1:** Take the middle of the  $t$ -space of reasearch:  $t_1 = 0.5$   
→ Highest values of  $F$  in  $[0, t_1]$

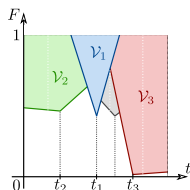
**Step 2:** Choose  $t_2$  in the middle of  $[0, t_1]$   
→ Highest values of  $F$  in  $[t_1, 1]$

**Step 3:** Choose  $t_3$  in the middle of  $[t_1, 1]$   
→ Highest values of  $F$  in  $[t_1, t_3]$

**Step 4:** Choose  $t_4$  in the middle of  $[t_1, t_3]$



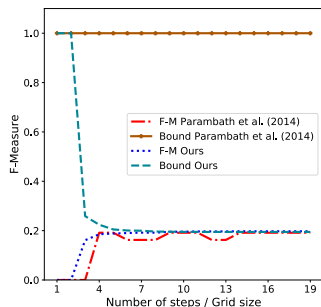
Step 3



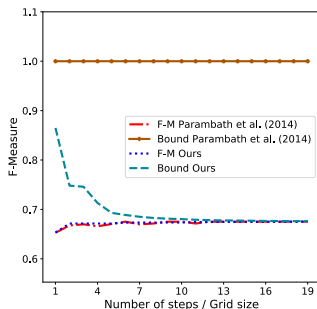
Step 4

# CONE: an Algorithm for F-measure Optimization

Comparison in terms of convergence



Abalone 12



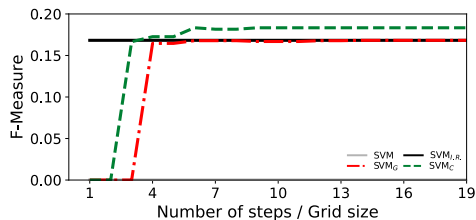
Adult

- A more informative bound
- A faster convergence
- Improves the performances

# CONE: an Algorithm for F-measure Optimization

## Comparison of performances

Abalone 12



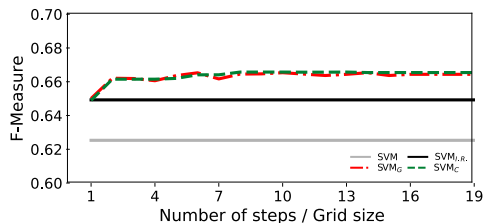
**SVM**: a linear SVM

**SVM<sub>I,R</sub>**: a linear SVM with weighted errors

**SVM<sub>G</sub>**: grid approach  
(Parambath et al., 2014)

**SVM<sub>C</sub>**: our approach

Adult



1. Introduction on Learning From Imbalanced Data
2. A Geometrical Approach based on the Distance to Positives
  - 2.1 Building Risky Areas

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*
  - 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*
3. An Approach based on Cost-Sensitive Learning
  - 3.1 Optimizing F-measure by Cost-Sensitive Classification

*CONE: "From Cost-Sensitive Classification to Tight F-Measure Bounds"*
  - 3.2 Improving the Benefits of Mass Distribution

*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*
4. Conclusion and Perspectives



# Improving Retailers Benefits

## Current model

### Currently

Model based on classification error (Decision Tree (Breiman et al., 1984) and Giny criterion)

Limits the number of alarms

Focuses on the number of false alarms, i.e. high precision

→ **Does not take main criterion into account: benefits**

# Improving Retailers Benefits

## Current model

### Currently

Model based on classification error (Decision Tree (Breiman et al., 1984) and Gini criterion)

Limits the number of alarms

Focuses on the number of false alarms, i.e. high precision

→ **Does not take main criterion into account: benefits**

### Idea

Define a new loss which optimizes retailers benefits

Use the amount in the loss function

# Improving Retailers Benefits

## Cost-Sensitive Model

Compute retailers benefits using a cost matrix (Elkan, 2001)

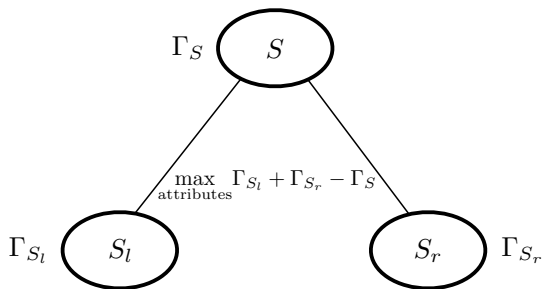
	Predicted Positive	Predicted Negative
Actual Positive	$C_{TP}$	$C_{FN}$
Actual Negative	$C_{FP}$	$C_{TN}$

$$\begin{aligned}C_{TP} &= 0 & C_{FN} &= (r - c(m)) \cdot m \\C_{FP} &= \rho \cdot r \cdot m - \xi & C_{TN} &= r \cdot m\end{aligned}$$

$$\ell(y, \hat{y}) = \sum_{i=1}^m [y_i(\hat{y}_i C_{TP_i} + (1 - \hat{y}_i) C_{FN_i}) + (1 - y_i)(\hat{y}_i C_{FP_i} + (1 - \hat{y}_i) C_{TN_i})].$$

# Improving Retailers Benefits

## Decision tree and splitting criterion 1/2



### Decision tree

Impurity:  $\Gamma = 1 - \sum_{i \in \mathcal{Y}} p_i^2$

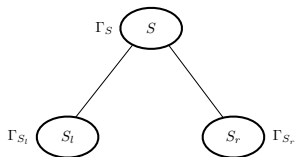
Split:  $\max_{attributes} \sum_{v \in \text{Children}} \Gamma_S - \alpha_v \Gamma_{S_v}$

### Weighted version

$$\Gamma_S \left[ \frac{1}{m} \sum_{i \in S_-} \left[ \frac{m_+}{m} c_{FP_i} + \frac{m_-}{m} c_{TN_i} \right] + \frac{1}{m} \sum_{i \in S_+} \left[ \frac{m_+}{m} c_{TP_i} + \frac{m_-}{m} c_{FN_i} \right] \right] =$$

# Improving Retailers Benefits

## Decision tree and splitting criterion 2/2

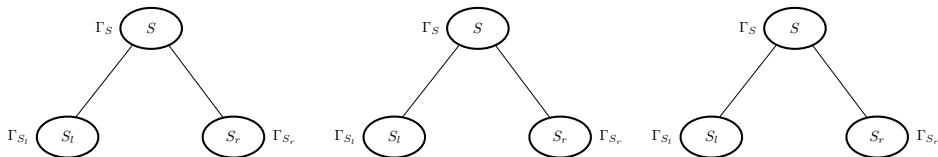


### Label

Choose the label that maximizes to profits

# Improving Retailers Benefits

## Decision tree and splitting criterion 2/2



### Label

Choose the label that maximizes to profits

### Random Forest

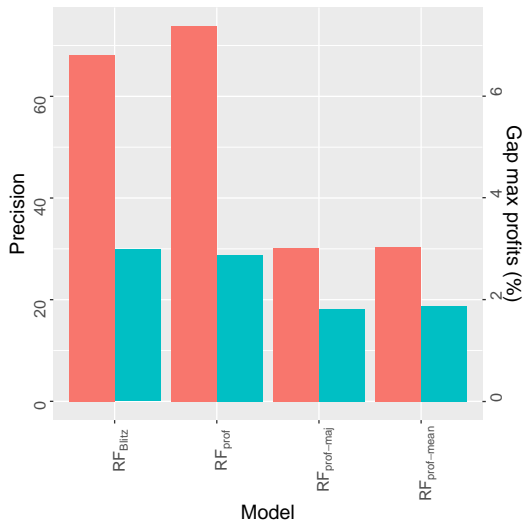
Build several decision trees using the splitting criterion  
Combination using different rules:

- simple majority vote

- weighted majority vote using the induced benefits

# Improving Retailers Benefits

## Experiments



4 months of transactions:

- Improves the profits
- Reduces the precision

A gap of 1% represents around 43 000 euros.

# Improving Retailers Benefits

## Gradient tree boosting

**Boosting:** Combine models such that  $f_t$  compensates for  $F_{t-1}$  weaknesses.

$$F_T = f_0 + \sum_{t=1}^T \alpha_t f_t$$

**Gradient Boosting:** Same idea, but work in the prediction space rather than parameter space.

$$g_t = - \left[ \frac{\partial \ell(y, F_{t-1}(\mathbf{x}_i))}{\partial F_{t-1}(\mathbf{x}_i)} \right], \quad (f_t, \alpha_t) = \underset{\alpha, f}{\operatorname{argmin}} \sum_{i=1}^m (r_i - \alpha f(x_i))^2.$$

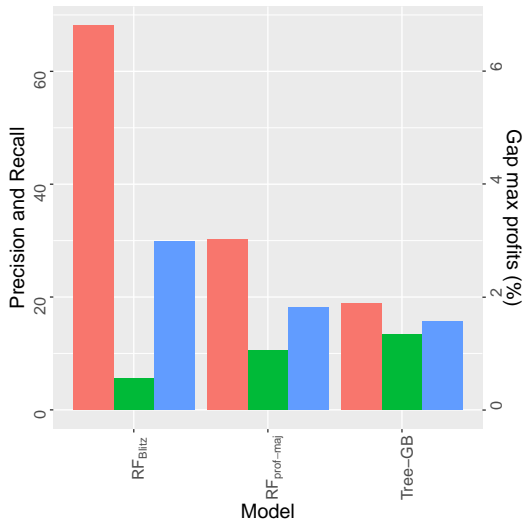
→ Give a surrogate of predefined  $\ell$  using the exponential

$$\ell(\mathbf{x}_i, y_i) = y_i(1 - c_i) \exp(-F(\mathbf{x}_i)) + c_i(1 - y_i) \exp(F(\mathbf{x}_i)).$$



# Improving Retailers Benefits

## Experiments



## Using Gradient Boosting

- Reduces training process
- Improves profits
- Higher recall
- Lower precision

Save around 60 000 euros compared to the current one

1. Introduction on Learning From Imbalanced Data
2. A Geometrical Approach based on the Distance to Positives
  - 2.1 Building Risky Areas

*ME<sup>2</sup> : "Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guaranties"*
  - 2.2 An Adjusted Version Nearest Neighbor Algorithm

*$\gamma - k$ -NN : "An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data"*
3. An Approach based on Cost-Sensitive Learning
  - 3.1 Optimizing F-measure by Cost-Sensitive Classification

*CONE: "From Cost-Sensitive Classification to Tight F-Measure Bounds"*
  - 3.2 Improving the Benefits of Mass Distribution

*"Tree-based Cost-Sensitive Methods for Fraud Detection in Imbalanced Data"*
4. Conclusion and Perspectives

# Conclusion

## Summary of Contributions

Two main axes were proposed to deal with the problem of learning from imbalanced data:

# Conclusion

## Summary of Contributions

Two main axes were proposed to deal with the problem of learning from imbalanced data:

1. Geometric: based on the distance to positives
  - Risky areas + local learning
  - Modification of the  $k$ -NN, modifying distance to positives

# Conclusion

## Summary of Contributions

Two main axes were proposed to deal with the problem of learning from imbalanced data:

### 1. Geometric: based on the distance to positives

- Risky areas + local learning
- Modification of the  $k$ -NN, modifying distance to positives

### 2. Cost-Sensitive Learning: Weighting the errors

- Bounds + iterative algorithm: optimizing F-measure
- Loss + algorithm: improving retailers benefits

# Conclusion

Bilan

	Advantages	Disadvantages
$ME^2$	Easy to learn $M$ Theoretical guarantees on $FP$	Over-fitting Detect new positives
$\gamma$ -k- NN	Easy to implement  Simplicity	Distance Computation  Too simple
CONE	Bounds on $F_\beta$ Derivation of an algorithm Require only few iterations	Algorithm convergence Guarantee at test time
$GB_{Tree}$	Fast to learn Flexibility	Low Precision

# Perspectives

## $\gamma$ -k-NN: a Metric Learning version

Based on the work on LMNN Weinberger and Saul (2009)

→ Propose a version of  $\gamma$ -k-NN based on learning new representations.

### Ideas :

- Keep compromise  $FN$  vs  $FP$ .
- Hyper-parameters: optimized to maximize the F-measure

### Deriving theoretical guarantees:

- On the learned metric (Bellet et al., 2015)
  - On the classification performances
- Ongoing work : submission at AISTATS 2020

# Perspectives

## CONE: Deriving lower bounds

**Lemma:** The difference  $(e_1 - e_2)$  is a decreasing function of  $t$  when  $\mathbf{e}(t)$  is obtained from an optimal classifier  $h$  learned with the weights  $\mathbf{a}(t)$ .



# Perspectives

## CONE: Deriving lower bounds

**Lemma:** The difference  $(e_1 - e_2)$  is a decreasing function of  $t$  when  $\mathbf{e}(t)$  is obtained from an optimal classifier  $h$  learned with the weights  $\mathbf{a}(t)$ .

Example when  $t' > t$ :

$$\begin{aligned} F(\mathbf{e}') - F(\mathbf{e}) &= \Phi_{\mathbf{e}} \left( \underbrace{\langle \mathbf{a}(t), \mathbf{e} \rangle}_{\text{blue}} + (t' - t)(e_2 - e_1) - \underbrace{\langle \mathbf{a}(t'), \mathbf{e}' \rangle}_{\text{red}} \right), \\ &= \Phi_{\mathbf{e}} \left( \underbrace{t(e_2 - e_1)}_{\text{green}} + \underbrace{(1 + \beta^2)e_1 - (1 + \beta^2)e'_1}_{\text{orange}} - t'(e'_2 - e'_1) + (t' - t)(e_2 - e_1) \right), \\ &\quad \downarrow \text{Use of the Lemma} \\ &\geq \Phi_{\mathbf{e}} \left( \underbrace{t(e'_2 - e'_1) - t'(e'_2 - e'_1)}_{\text{blue}} + (1 + \beta^2)(e_1 - e'_1) + (t' - t)(e_2 - e_1) \right), \\ &\quad \downarrow \dots \\ F(\mathbf{e}') - F(\mathbf{e}) &\geq \Phi_{\mathbf{e}} \left( (1 + \beta^2)(e_1 - e'_1) + (t' - t)e_2 - e_1 - (e'_2 - e'_1) \right). \end{aligned}$$

# Perspectives

## CONE: Deriving lower bounds

$$F(\mathbf{e}') - F(\mathbf{e}) \geq \Phi_{\mathbf{e}} \left( (1 + \beta^2)(\mathbf{e}_1 - \mathbf{e}'_1) + (t' - t)\mathbf{e}_2 - \mathbf{e}_1 - (\mathbf{e}'_2 - \mathbf{e}'_1) \right)$$

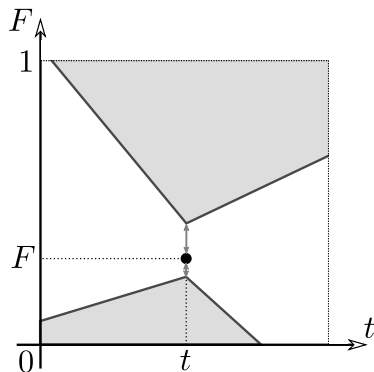
- $\mathbf{e}'_2 - \mathbf{e}'_1$ : as seen previously.
- $\mathbf{e}_1 - \mathbf{e}'_1$ : find a tight lower-bound.

# Perspectives

## CONE: Deriving lower bounds

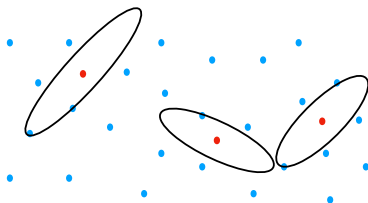
$$F(\mathbf{e}') - F(\mathbf{e}) \geq \Phi_{\mathbf{e}} \left( (1 + \beta^2)(e_1 - e'_1) + (t' - t)e_2 - e_1 - (e'_2 - e'_1) \right)$$

- $e'_2 - e'_1$ : as seen previously.
- $e_1 - e'_1$ : find a tight lower-bound.



- Bound the values of  $F_{\beta}$
- Get a new algorithm
- Deriving generalization bounds  $F_{\beta}$
- Optimality of  $F_{\beta}$  at test time
- Empirically : generalization bounds based on the **validation** (Kawaguchi et al., 2017).

**Problem:**  $ME^2$  is prone to over-fitting



→ Find a way to "smooth" the classification process

→ Convex Combinations of local models (Zantedeschi et al., 2016)

# Thank you for your attention !

## International Journal

- G.Metzler, X.Badiche, B.Belkasmı, E.Fromont, A.Habıard and M.Sebban; *Learning Maximum Excluding Ellipsoids from Imbalanced Data with Theoretical Guarantees*, PRL, 2018.

## International Conferences

- R.Viola, R.Emonet, A.Habıard, G.Metzler, S. Riou and M.Sebban; *An Adjusted Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data*, ICTAI, 2019.
- K.Bascol, R.Emonet, E.Fromont, A.Habıard, G.Metzler and M.Sebban; *From Cost-Sensitive Classification to Tight F-Measure Bounds*, AISTATS, 2019.
- G.Metzler, X.Badiche, B.Belkasmı, E.Fromont, A.Habıard and M.Sebban; *Tree-based Cost Sensitive Methods for Fraud Detection in Imbalanced Data*, IDA, 2018.

## National Conferences

- R.Viola, R.Emonet, A.Habıard, G.Metzler, S.Riou and M.Sebban; *Une version corrigée de l'algorithme des plus proches voisins pour l'optimisation de la F-mesure dans un contexte déséquilibré*, CAp, 2019.
- K.Bascol, R.Emonet, E.Fromont, A.Habıard, G.Metzler and M.Sebban; *Un algorithme d'optimisation de la F-Mesure par pondération des erreurs de classification*, CAp, 2018.
- G.Metzler, X.Badiche, B.Belkasmı, E.Fromont, A.Habıard and M.Sebban; *Apprentissage de Sphères Maximales d'exclusion avec Garanties Théoriques*, CAp, 2017.

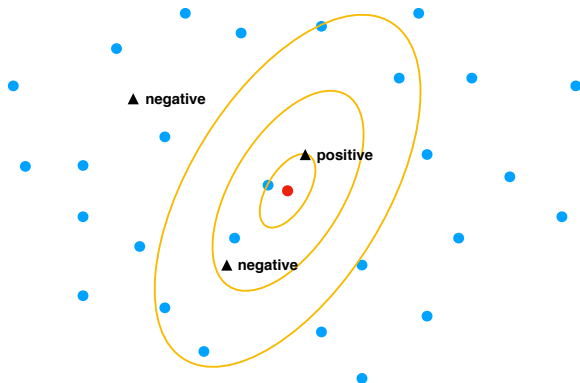
# References I

- Bellet, A., Habrard, A., and Sebban, M. (2015). Metric learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 9(1):1–151.
- Bousquet, O. and Elisseeff, A. (2002). Stability and generalization. *Journal of Machine Learning Research*, 2:499–526.
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and regression trees*. The Wadsworth statistics/probability series. Wadsworth and Brooks/Cole Advanced Books and Software, Monterey, CA.
- Elkan, C. (2001). The foundations of cost-sensitive learning. In *Proceedings of the 17th International Joint Conference on Artificial Intelligence - Volume 2*, pages 973–978, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Kawaguchi, K., Kaelbling, L. P., and Bengio, Y. (2017). Generalization in deep learning. *arXiv preprint arXiv:1710.05468*.
- Parambath, S. P., Usunier, N., and Grandvalet, Y. (2014). Optimizing f-measures by cost-sensitive classification. In *Advances in Neural Information Processing Systems (NIPS-14)*, pages 2123–2131.
- Tax, D. M. J. and Duin, R. P. W. (2004). Support vector data description. *Machine Learning Journal*, 54(1):45–66.
- Weinberger, K. Q. and Saul, L. K. (2009). Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244.
- Zantedeschi, V., Emonet, R., and Sebban, M. (2016). Metric learning as convex combinations of local models with generalization guarantees. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

# $ME^2$ : Learning Risky Areas

## Algorithm

1. assign each text example to its closest positive
2. apply the following classification rule:



# Perspectives

## $\gamma$ -k-NN: Version Metric Learning

Based on the work on LMNN Weinberger and Saul (2009)

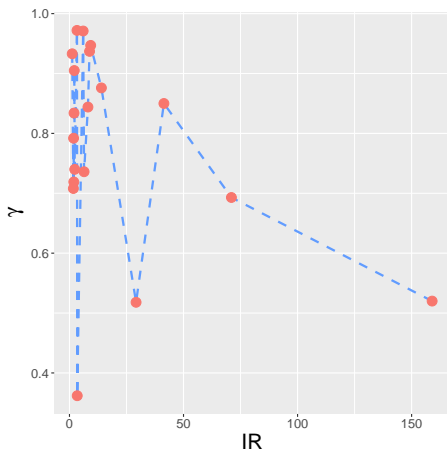
→ Propose a version of  $\gamma$ -k-NN based on learning new representations.

$$\min_{\mathbf{M} \in \mathcal{S}^+} \frac{1}{m^3} \left( \frac{1-\alpha}{2} \sum_{\substack{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{S} \\ y_i = y_j = 1}} d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j)^2 + \right. \\ \left. \frac{1-\alpha}{2} \sum_{\substack{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in \mathcal{R} \\ y_i = 1}} [1 - m' + d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j)^2 - d(\mathbf{x}_i, \mathbf{x}_k)]_+ \right. \\ \left. + \alpha \sum_{\substack{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in \mathcal{R} \\ y_i = -1}} [1 - m' + d(\mathbf{x}_i, \mathbf{x}_j)^2 - d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_k)]_+ \right) \mu \|\mathbf{M} - \mathbf{I}\|_{\mathcal{F}}^2.$$



# $\gamma$ -k-NN: a revisit of the $k$ -NN

$\gamma^*$  vs. I.R.

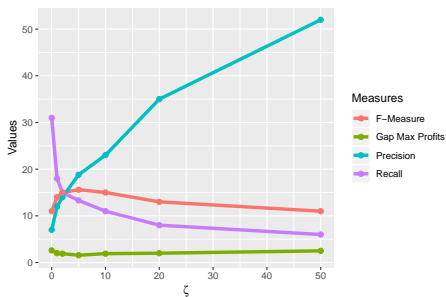


→ In average  $\gamma$  is a decreasing function of  $IR$

# Improving Retailers Benefits

## Study of parameter $\xi$

### Influence of $\xi$ parameter



### Increasing $\xi$ value:

Improves the Precision

Reduces the Recall

Reduces the retailers benefits

# Comparison on Blitz dataset

## Comparison Contributions

